

## Article

# *In Silico* Prediction of siRNA Ionizable-Lipid Nanoparticles *in vivo* Efficacy: Machine Learning Modeling Based on Formulation and Molecular Descriptors

Abdelkader A Metwally <sup>1,2,\*</sup>, Amira A Nayel <sup>3</sup> and Rania M Hathout <sup>2</sup>

<sup>1</sup> Department of Pharmaceutics, Faculty of Pharmacy, HSC, Kuwait University, Kuwait; abdelkader.metwally@ku.edu.kw (A.A.M.)

<sup>2</sup> Department of Pharmaceutics and Industrial Pharmacy, Faculty of Pharmacy, Ain Shams University, Egypt; abdelkader\_ali@pharma.asu.edu.eg (A.A.M.); rania.hathout@pharma.asu.edu.eg (R.M.H.)

<sup>3</sup> Clinical Pharmacy Department, Alexandria Ophthalmology Hospital, Alexandria, Egypt; amiranayel@gmail.com (A.A.N)

\* Correspondence: abdelkader.metwally@ku.edu.kw

**Abstract:** *In silico* prediction of the *in vivo* efficacy of siRNA ionizable-lipid nanoparticles is desirable yet never achieved before. This study aims to computationally predict siRNA nanoparticles *in vivo* efficacy, which saves time and resources. A data set containing 120 entries was prepared by combining molecular descriptors of the ionizable lipids together with two nanoparticles formulation characteristics. Input descriptor combinations were selected by an evolutionary algorithm. Artificial neural networks, support vector machines and partial least squares regression were used for QSAR modeling. Depending on how the data set is split, two training sets and two external validation sets were prepared. Training and validation sets contained 90 and 30 entries respectively. The results showed the successful predictions of validation set log(dose) with  $R^2_{val} = 0.86 - 0.89$  and  $0.75 - 80$  for validation sets one and two respectively. Artificial neural networks resulted in the best  $R^2_{val}$  for both validation sets. For predictions that have high bias, improvement of  $R^2_{val}$  from 0.47 to 0.96 was achieved by selecting the training set lipids lying within the applicability domain. In conclusion, *in vivo* performance of siRNA nanoparticles was successfully predicted by combining cheminformatics with machine learning techniques.

**Keywords:** siRNA; ionizable lipids; nanoparticles; *in vivo*; QSAR; machine learning

## 1. Introduction

The process of developing short interfering RNA (siRNA) lipid nanoparticles is lengthy and time consuming because it involves the initial chemical synthesis of a usually large number of ionizable lipids and lipid-like molecules [1-3], the formulation of siRNA nanoparticles and the subsequent *in vitro* and *in vivo* evaluation of these nanoparticles, in an attempt to find the best ionizable lipid that is suitable for clinical use in terms of efficacy and safety. Alnylam's small interfering RNA (siRNA) stable nucleic acid lipid nanoparticles, currently marketed as Onpattro™ (Patisiran), obtained FDA approval in 2018. This was followed by FDA approval of Alnylam's Givosiran™ and Lumasiran™ in 2019 and 2020 respectively [4].

Gene silencing by double-stranded RNA (dsRNA) was reported by Fire and Mello in *Caenorhabditis elegans* [5] and later siRNA duplexes of length 21-22 nucleotides proved to promote post-transcriptional gene silencing in mammalian cells [6]. Since then, the potential of siRNA as a therapeutic macromolecule against many diseases has been investigated, with more than 40 siRNA based therapies already reaching phases 2, 3 or 4 of clinical trials.[7-9] The major barriers against the successful employment of therapeutic siRNA include the lack of stability of the siRNA duplex, the immune response mediated

by TOL-like receptors, the rapid renal clearance of naked siRNA, and the difficulty of the intracellular delivery of unmodified siRNA due to its large size and the large number of negative charges on its back-bone [10,11].

One method to overcome the barriers of siRNA delivery is to formulate it as siRNA ionizable lipid nanocomplexes (lipoplexes) or lipidic nanoparticles [12-15]. These nanoparticles are multicomponent and may also contain helper lipids, PEG-lipids and phospholipids. An ideal delivery system should ensure response reproducibility, non-immunogenicity, good payload and ease of manufacturing [13].

The process of preparing siRNA lipoplexes and nanoparticles involves many steps: the synthesis of the ionizable lipids, their purification and characterization, then the process of preparing the nanoparticles including determining the siRNA to cationic lipid ratio, the cationic lipid to helper lipid (if any) ratio, and nanoparticles characterization in terms of their size, zeta potential, pK<sub>a</sub>, stability and *in vivo* evaluation of their safety and silencing efficacy. All of these steps require time and resources and indeed if the *in vivo* efficacy, as measured by either the siRNA dose or knockdown efficiency, could be predicted within reasonable accuracy by using computational means, the process of developing siRNA nanomedicines would be vastly improved in terms of time and costs. Therefore it is important to attempt to predict the *in vivo* efficacy of siRNA cationic lipid nanoparticles by using machine learning techniques. These techniques can be generally classified into two main groups; supervised and unsupervised learning methods. Supervised learning is used in tasks such as regression and classification, i.e. when there is a dependent variable and one or more independent variables.

Artificial neural networks (ANNs) are a collection of linear and non-linear functions that map an input to an output. These functions can approximate a nonlinear complex function. The idea behind the inner working of ANNs is that input data ( $\mathbf{x}$ ) are scaled and combined in a linear manner in the form of  $\mathbf{W}\mathbf{x} + \mathbf{b}$ , where  $\mathbf{W}$  is the weights matrix and  $\mathbf{b}$  is bias, and then the output of this linear combination is fed into a non-linear function (called activation function), the output of which could be used as an input to the next layer and/or to a final output layer [16].

Support vector machines (SVM) are a supervised machine learning technique. For classification, SVM aims to find a hyperplane (decision surface) that can separate two classes of observations with a maximum margin of separation [17]. Similarly, SVM regression follows the same logic of finding a hyperplane, but with a fixed margin width, epsilon ( $\epsilon$ ), within which the prediction error is considered zero, and the hyperplane found should minimize the sum of squared error, i.e., the sum of the difference between the actual and predicted values:  $\sum_{i=1}^n y_i - (\mathbf{W}^T \mathbf{x}_i + \mathbf{b})$ . To enable the formulation of non-linear decision surfaces, a kernel function is applied. The general form of the kernel functions is  $K(\mathbf{x}_1, \mathbf{x}_2) = \langle \varphi(\mathbf{x}_1), \varphi(\mathbf{x}_2) \rangle$ , where  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are two data points. The kernel function thus avoids the actual calculation of the function  $\varphi$  [18].

Partial least squares (PLS) regression is another supervised learning technique [19]. PLS combines dimensionality reduction of the data with a regression model. PLS formulation of the latent variables (scores or components) is carried out with the aim of maximizing the covariance of the components with the response variable, which differentiates PLS from regular principle component analysis (PCA) [20]. The response variable in PLS may be univariate or multivariate. For the prediction of a new data point response  $\hat{y}'_o$  from a predictor point  $x'_o$ , the following equation applies:  $\hat{y}'_o = \frac{1}{n} \sum_{i=1}^n y'_i + \mathbf{B}^T (x_o - \frac{1}{n} \sum_{i=1}^n x'_i)$ .  $\mathbf{B}$  is the matrix of regression coefficients, and is defined as:  $\mathbf{B} = \mathbf{W}(\mathbf{T}^T \mathbf{T})^{-1} \mathbf{T}^T \mathbf{Y}$ , where  $\mathbf{W}$  is the matrix of weights and  $\mathbf{T} = \mathbf{XW}$  [20].

In order to extract chemical information from the structures of the molecules under investigation, molecular descriptors, which are important cheminformatics tools, are employed to carry out this task [21,22]. Molecular descriptors are numerical values resulting from either an experimental procedure or from a set of mathematical and/or logical algorithms that are performed on chemical structures [23]. The descriptors can be generally classified as 0D and 1D, when only molecular formula or constitutional properties of a molecule are considered, while 2D descriptors are calculated based on topo-

logical properties of a molecule and 3D descriptors depend on geometrical properties of a molecule. Further classifications include 2.5D chiral descriptors and descriptors with more than three dimensions [24,25]. Molecular descriptors have been used as predictors of the self-assembly of drug molecules into nanoparticles [26], to model drug binding kinetics [27], in QSAR modeling [28] and in target identification [29]. Molecular descriptors were also used to successfully predict the binding energy between drug molecules and their nanocarriers and hence predict drug loading onto lipidic and polymeric nanoparticles [30].

Previous QSAR studies on nanoparticles have mostly addressed predicting the cellular uptake and toxicological properties of inorganic nanoparticles, with either unmodified or modified surfaces [31-33], however, developing QSAR models for predicting siRNA *in vivo* efficacy has not been achieved before.

In the current work, a data set is prepared using five publications [1,34-37]. This data set contains the 1D and 2D descriptors of ionizable lipids together with both of the formulation descriptor (PEG mol%) and the percentage knockdown resulting from a specific siRNA dose. The siRNA nanoparticles *in vivo* efficacy when formulated with these ionizable lipids was included as the response variable; logarithm of the dose resulting in a specific knockdown percent of the target gene. The data set is split into training and validation sets, where the training set is used to construct the machine learning models, and the validation set is used as an external test set that is used only to evaluate the predictive models constructed by modeling the training set. An evolutionary algorithm is used to select the best descriptor combinations and is combined with three machine learning techniques; ANN, SVM and PLS regression, to build the predictive models. The performance of the predictive models using the three machine learning techniques and the quality of predictions and how to improve them is presented and discussed. Figure 1 shows the work flow of the modeling and evaluation process.

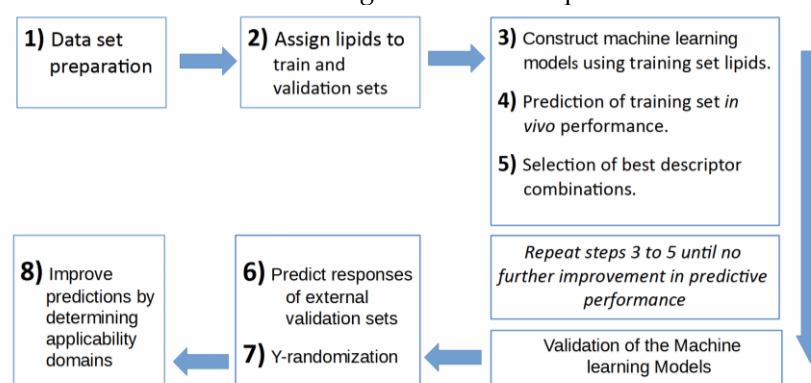


Figure 1. The workflow of the predictive model building process.

## 2. Materials and Methods

### 2.1 Data set preparation

#### 2.1.1 Data selection from available literature

For preparing the data set, five publications[1,34-37] were retrieved after carrying out online search using PUBMED and Google Scholar servers, where all of them fulfilled the following requirements: siRNA is delivered by means of ionizable lipids, siRNA *in vivo* performance is evaluated *in vivo* against factor FVII expression, all nanoparticles contained the ionizable lipid, DSPC, cholesterol and PEG-lipid (with PEG average molecular weight = 2000), and the PEG-lipid mole % in the formulation is either given or can be calculated. In addition, both the siRNA dose and the percentage knockdown or percentage gene expression resulting from a specific siRNA dose must be provided. Five papers were selected to prepare the data set [1,34-37]. Wherever the values for the gene

expression or dose were not provided numerically, these values were obtained from the relative figures using WebPlotDigitizer v4.2. In case two or more lipids had the same 2D structure, one of them was retained. If an ionizable lipid lacked a well defined *in vivo* efficacy measure, such as a definite dose or knockdown %, it was omitted.

### 2.1.2 Calculation of the 2D molecular descriptors

The structures of the ionizable lipids were drawn using ACD ChemsSketch, and the structures were saved as either individual MDL .mol files or combined together into a single .sdf file using OpenBabel v2.4 [38]. The following software packages were used for the calculation of the 1D/2D molecular descriptors: Padel Descriptor v2.21 [39], RDKit 2017, and ToMoCoMD QuBiLS-MAS 2020 [25]. For the calculation of the QuBiLS-MAS descriptors, the following settings were selected: linear algebraic form, atom-based, non-stochastic matrix form, and total groups.

### 2.1.3 Data set preprocessing

The initial data set containing the descriptors was further processed by removing columns having one or more of either missing or not available (NA) entries. Columns with same-value entries were also removed. If certain columns in the data set showed a high correlation (cutoff  $r = 0.98$ ) with each other [40], all the columns were removed except for one column which has the lowest average correlation with the other descriptor (predictor) columns in the data set. In addition, the formulation descriptor (PEG mol%) and percentage knockdown resulting from a specific siRNA dose were added as predictors. The data set descriptor columns were scaled by calculating the z-scores. The siRNA nanoparticles *in vivo* efficacy was included as the response variable; logarithm of the dose resulting in a specific knockdown percent.

## 2.2. Principle component analysis (PCA) of data set

PCA of the scaled data set predictor columns (without response columns) was carried out using ChemometricsWithR package through R software v3.5.

## 2.3. Splitting the data set into training and validation sets

For modeling purposes, the data set entries were split into a training set (75% of entries) and a validation set (25% of entries). This process was carried out two times separately on the data set where the validation set entries (or observations) were selected either by random selection or by selecting sequentially every fourth entry in the set, with the remainder of the entries in the data set taken as the training set.

## 2.4. Machine learning models

The modeling process was carried out using either R software version 3.5 or Microsoft Open R v3.5. The following R packages were used for all modeling methods: caret [41] and Metrics [42]. For artificial neural networks modeling, nnet package was used. The hyperparameters were one hidden layer, two nodes and a weight decay of 0.1 for training and 0.001 for final validation set predictions. The support vector machine regression modeling (epsilon-regression) was carried out using kernlab package [43], with epsilon value of 0.1 and the kernel chosen to be the Gaussian radial basis function kernel defined as  $K(\mathbf{x}, \mathbf{x}_i) = -\sigma ||\mathbf{x} - \mathbf{x}_i||^2$ , where  $\sigma$  is the inverse width parameter and is determined by the package's sigest function. The partial least squares modeling was carried out using pls package [44] with the number of principle components covering 98% of the variance.

## 2.5. Selection of the molecular descriptors by the evolutionary algorithm

An evolutionary algorithm was written as an R script to select the best descriptors for model building. 400 initial parent combinations of descriptors were randomly selected, and then each one of them was used as an input to construct the machine learning models that are used to predict the training set log(dose) values and their associated RMSEs (training RMSE).

The training RMSE is calculated as follows: the training set is split into three folds, two folds are used to construct the machine learning model, and the third fold is used as a test set to calculate training RMSE. After evaluating the training RMSE for all predictor combinations, the best combinations are kept as parents and are used to construct offspring combinations. The process is repeated until no further improvement in training RMSE for this specific test fold. The whole selection process is repeated for each of the remaining two test folds. The parameters for the evolutionary algorithm are as follows: population size 400, 25% elitism, 20% mutation, number of generations 10-20 and multipoint cross-over.

RMSE is calculated as:

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (P_i - A_i)^2}{n}}$$

Bias is calculated as:

$$\text{Bias} = P_i - A_i$$

Where  $P_i$  and  $A_i$  are the predicted and actual log(dose) values of observation (lipid or entry)  $i$  respectively, and  $n$  is the number of observations.

## 2.6. Ensemble learning by averaging of the validation set predictions

The best descriptor combinations that result in the lowest training RMSE were used as inputs for the machine learning modeling algorithm that was used in the training; either ANN, SVM or PLS regression. The central tendency of the validation set predictions were calculated as median of these values for each validation set lipid. The validation set RMSE ( $\text{RMSE}_{\text{val}}$ ) and coefficient of determination ( $R^2_{\text{val}}$ ) were calculated using these median values. The  $R^2_{\text{val}}$  is calculated as:

$$R^2_{\text{val}} = \frac{(\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}))^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}$$

where  $x_i$  and  $y_i$  are the  $i^{\text{th}}$  predicted (the median value) and actual responses respectively,  $\bar{x}$  and  $\bar{y}$  are the mean values of predicted and actual responses respectively.

## 2.7. Y-Randomization of data set

To evaluate the validity of the resulting descriptor combinations, and the possibility that the obtained validation set predictions might be due to random chance, a Y-randomization of the training data set was carried out by randomizing the training set responses [45]. The predictive models were then constructed by using these randomized responses for model training and subsequent validation as described in section 2.6.

## 3. Results

### 3.1. Data set preprocessing and preparation

The number of observations included in the data set after omitting the lipids or entries that fit the omitting criteria explained in section 2.1.1 was 120 entries (rows). The resulting data set contained 438 predictor columns; 436 columns of molecular de-



scriptors, and 2 columns for PEG mol % and knockdown %. In addition, one response column was included; logarithm of siRNA dose that results in a specific knockdown of the target gene. Table 1 provides summary of the data set.

**Table 1.** Summary of data set. The entries represent either distinct lipids or the same lipid but with different PEG mol % and/or knockdown %.

Index of entries	Number of entries per study	Reference
1-30	30	K. Rajappan et al.[34]
31-62	32	C. A. Alabi et al.[37]
63-95, 105	34	M. Jayaraman et al.[1]
96-104	9	V. Kumar et al.[36]
106-120	15	K. A. Whitehead et al.[35]

### 3.2. Splitting the data set into training and validation sets

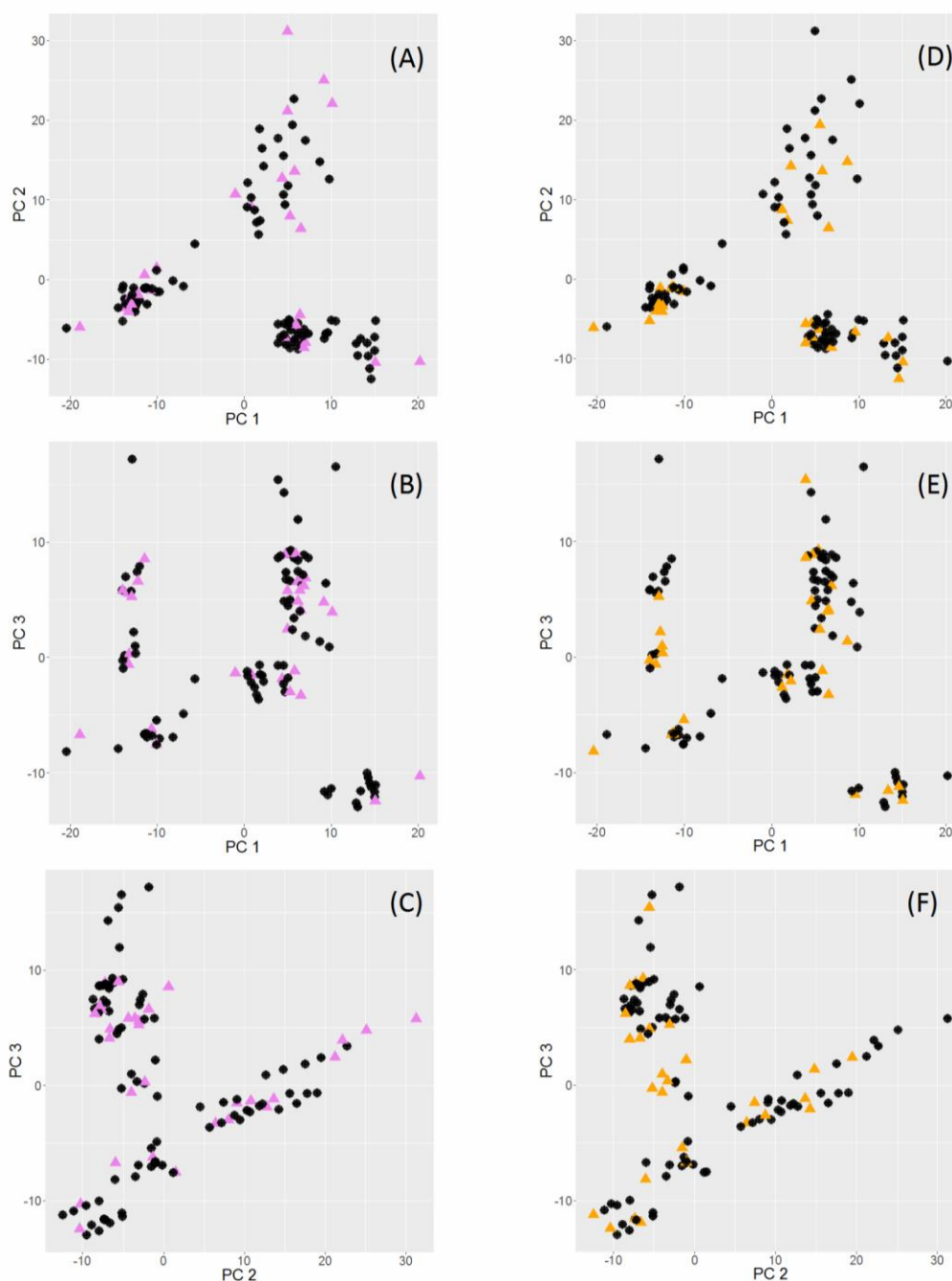
Two different methods were used to select the validation set entries, with the remainder of the entries in each splitting method being used for training the machine learning models. These selection processes resulted in the following data sets: training set1, validation set 1, training set 2 and validation set 2. These sets are shown in Table 2. Each training and validation set contained 90 and 30 entries, respectively.

To investigate the relationships between the training and validation entries, PCA was carried out, to project the data points on the newly formed principle components, capturing as much as possible of the variance of the data. The PCA score plots are shown in Figure 2. Principle components 1, 2 and 3 (PC 1, PC 2 and PC 3) contributed to 22%, 19% and 12% of the total variance, respectively. The observations of validation set 1 and 2, shown as colored triangles, show homogenous spread among those of training set 1 and 2 respectively, which is an important characteristic of any training and external validation samples, as the training set must reasonably represent the characteristics of the validation set as well as capturing the general characteristics of the whole data set. Both splitting methods of the data set, whether random splitting or sequential selection of the validation entries, resulted in good spread of the validation entries among the training ones, with no significant presence of outlier observations of the validation sets with respect to their respective training set.

Table 2. Training and validation sets 1 and 2.

Set	Training entries index	Validation entries index
<b>1</b>	3, 4, 5, 6, 8, 9, 10, 11, 13, 14, 18, 20, 21, 24, 25, 26, 27, 28, 29, 30, 33, 34, 35, 36, 37, 39, 40, 41, 42, 43, 45, 46, 47, 48, 49, 51, 54, 55, 56, 57, 58, 59, 60, 62, 63, 64, 66, 67, 68, 69,	1, 2, 7, 12, 15, 16, 17, 19, 22, 23, 31, 32, 38, 44, 50, 52, 53, 61, 65, 70, 74, 76, 77, 81, 85, 87, 88, 91, 109, 116.

	71, 72, 73, 75, 78, 79, 80, 82, 83, 84, 86, 89, 90, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 110, 111, 112, 113, 114, 115, 117, 118, 119, 120.	
2	1, 2, 3, 5, 6, 7, 9, 10, 11, 13, 14, 15, 17, 18, 19, 21, 22, 23, 25, 26, 27, 29, 30, 31, 33, 34, 35, 37, 38, 39, 41, 42, 43, 45, 46, 47, 49, 50, 51, 53, 54, 55, 57, 58, 59, 61, 62, 63, 65, 66, 67, 69, 70, 71, 73, 74, 75, 77, 78, 79, 81, 82, 83, 85, 86, 87, 89, 90, 91, 93, 94, 95, 97, 98, 99, 101, 102, 103, 105, 106, 107, 109, 110, 111, 113, 114, 115, 117, 118, 119.	4, 8, 12, 16, 20, 24, 28, 32, 36, 40, 44, 48, 52, 56, 60, 64, 68, 72, 76, 80, 84, 88, 92, 96, 100, 104, 108, 112, 116, 120.



**Figure 2.** PCA score plots. A - C: training and validation set 1 entries are shown as black circles and pink triangles respectively. D - F: training and validation set 2 entries are shown as black circles and orange triangles respectively.

### 3.3. Selection of the molecular descriptors by the evolutionary algorithm

When constructing the descriptor combinations to be used as inputs for the machine learning algorithm, the PEG mol % and the knockdown % were always included in the combinations. Any additional molecular descriptors were added and selected by the evolutionary algorithm. Figure 3 shows the top six molecular descriptors with the highest frequencies of appearance in the descriptor combinations that are selected by the evolutionary algorithm. For each machine learning method, ANN, SVM or PLS, the descriptor with highest frequency was considered 100 % and the other descriptors frequencies were calculated relative to it. It is evident that each machine learning model resulted in different top descriptors. It is also clear that the training sets one and two resulted in different top descriptors for the same machine learning method. The only common descriptors, taking the two training sets and the three machine learning meth-



ods in consideration, were PEOE\_VSA9, GATS3m and GATS8p. PEOE\_VSA9 is a Van der Waals surface area descriptor that describes atomic partial charges. GATS3m and GATS8p are Geary autocorrelation - lag 3 weighted by atomic masses and Geary autocorrelation - lag 8 weighted by atomic polarizabilities respectively. It should be noted that these descriptors are present in combinations of descriptors (predictors) including the PEG mol % and the knockdown %, thus, their direct influence on the *in vivo* performance of the ionizable lipids should be limited to this context.

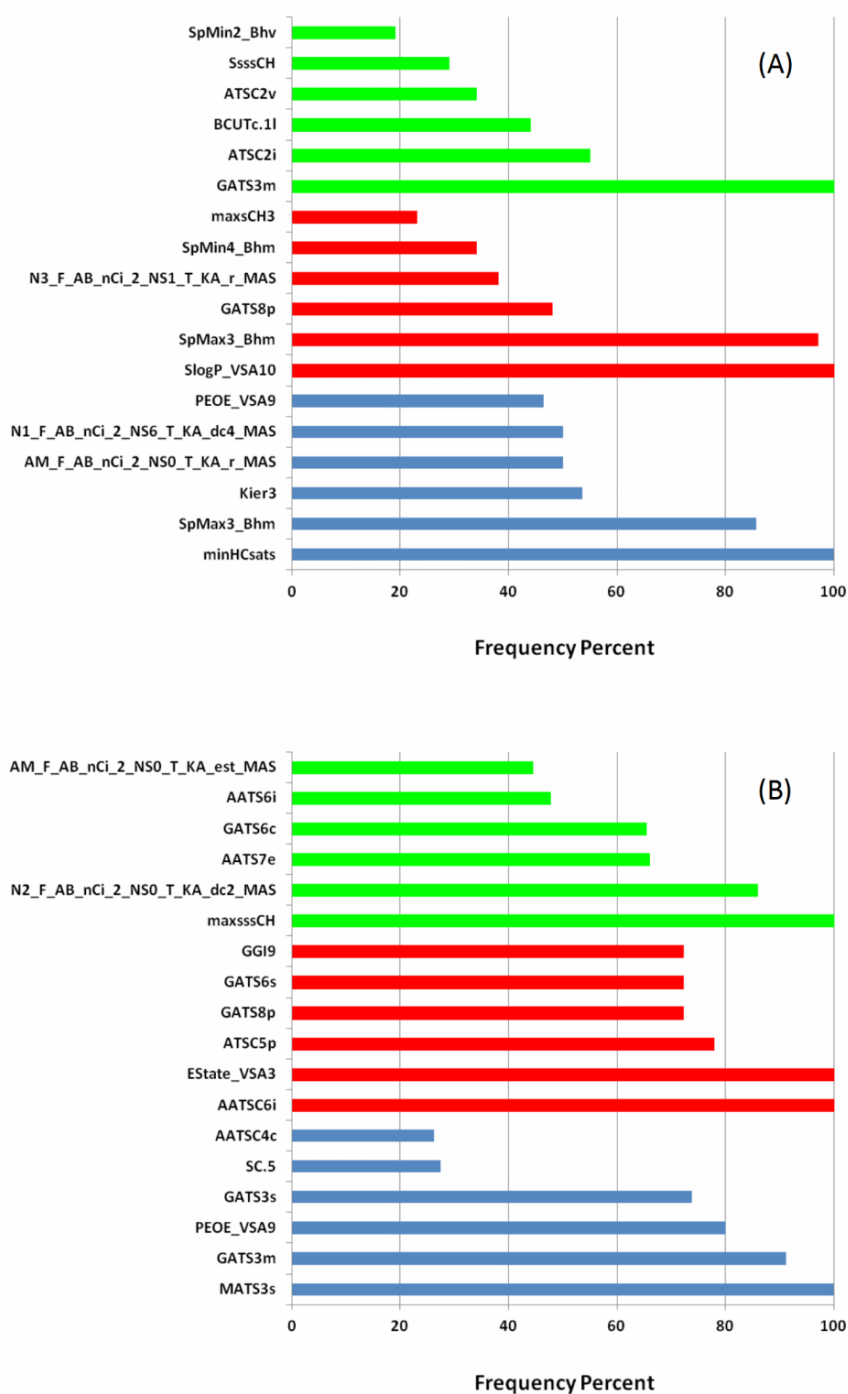
The number of molecular descriptors in each descriptor combination as selected by the evolutionary algorithm is listed in Table 3. It is to be noted that these molecular descriptors are present in each combination in addition to both PEG mol % and knockdown %, with the later two being present in each predictor combination. The number of final combinations for all methods for each training set was 300 combinations. It was noticed that there were repeated combinations in the final 300 combinations, as omission of descriptors by the evolutionary algorithm results eventually in offspring combinations of the same descriptors. For example, there were 73 unique combinations among the final 300 combinations selected by the evolutionary algorithm and ANN training of set one.

**Table 3.** The minimum, maximum and median number of the molecular descriptors in the final predictor combinations for each training set and machine learning method.

Training Set	Machine learning method	min	max	median
1	ANN	2	7	5
1	SVM	3	7	4
1	PLS	3	7	3
2	ANN	4	9	5
2	SVM	4	9	6
2	PLS	4	9	6

The improvement in predictions of the validation set responses at the end of the evolutionary algorithm is shown in Table 4. The  $RMSE_{val}$  in the table are calculated as the first quartile of the RMSE of predictions using the initial 400 descriptor combinations and the final 400 descriptor combinations at the end of the evolutionary algorithm iterations. It is clear that there were improvement in the quality of individual predictions for both validation sets and for all methods as evident by the decrease in the  $RMSE_{val}$ .

The predictive performance of the machine learning models was evaluated by predicting the validation sets responses. The validation sets were neither used in the selection of best descriptor combinations by the evolutionary algorithm nor they were used in the training of the predictive models, thus, the validation sets represent external unknown test samples for the machine learning models. Using the descriptor combinations selected by the evolutionary algorithm, the median (averaged) predictions of the validation sets one and two resulted in  $R^2_{val}$  of 0.72 to 0.89 and  $RMSE_{val}$  of 0.23 to 0.36 (Table 5). The machine learning method used to predict the validation set responses had a strong effect on the predictive performance, with the ANN predictions resulting in the highest  $R^2_{val}$  of 0.89 and 0.80 for validation sets one and two respectively. Similarly, ANN resulted in the lowest  $RMSE_{val}$  of 0.23 and 0.30 for validation sets one and two respectively. There were also a difference in the predictive performance between validation sets one and two (Table 5), which reflects the effect of both the training set and validation sets compositions.



**Figure 3.** Relative frequencies of descriptors in the descriptor combinations selected by the evolutionary algorithm. A: training set 1. B: training set 2. Blue: ANN, red: SVM and green: PLS.

**Table 4.** Improvement of quality of individual validation set predictions by the evolutionary algorithm.

Validation Set	Machine learning method	Initial First quartile RMSE <sub>val</sub>	Final First quartile RMSE <sub>val</sub>
----------------	-------------------------	--	--

1	ANN	0.41	0.33
1	SVM	0.40	0.31
1	PLS	0.41	0.29
2	ANN	0.40	0.35
2	SVM	0.39	0.36
2	PLS	0.44	0.37

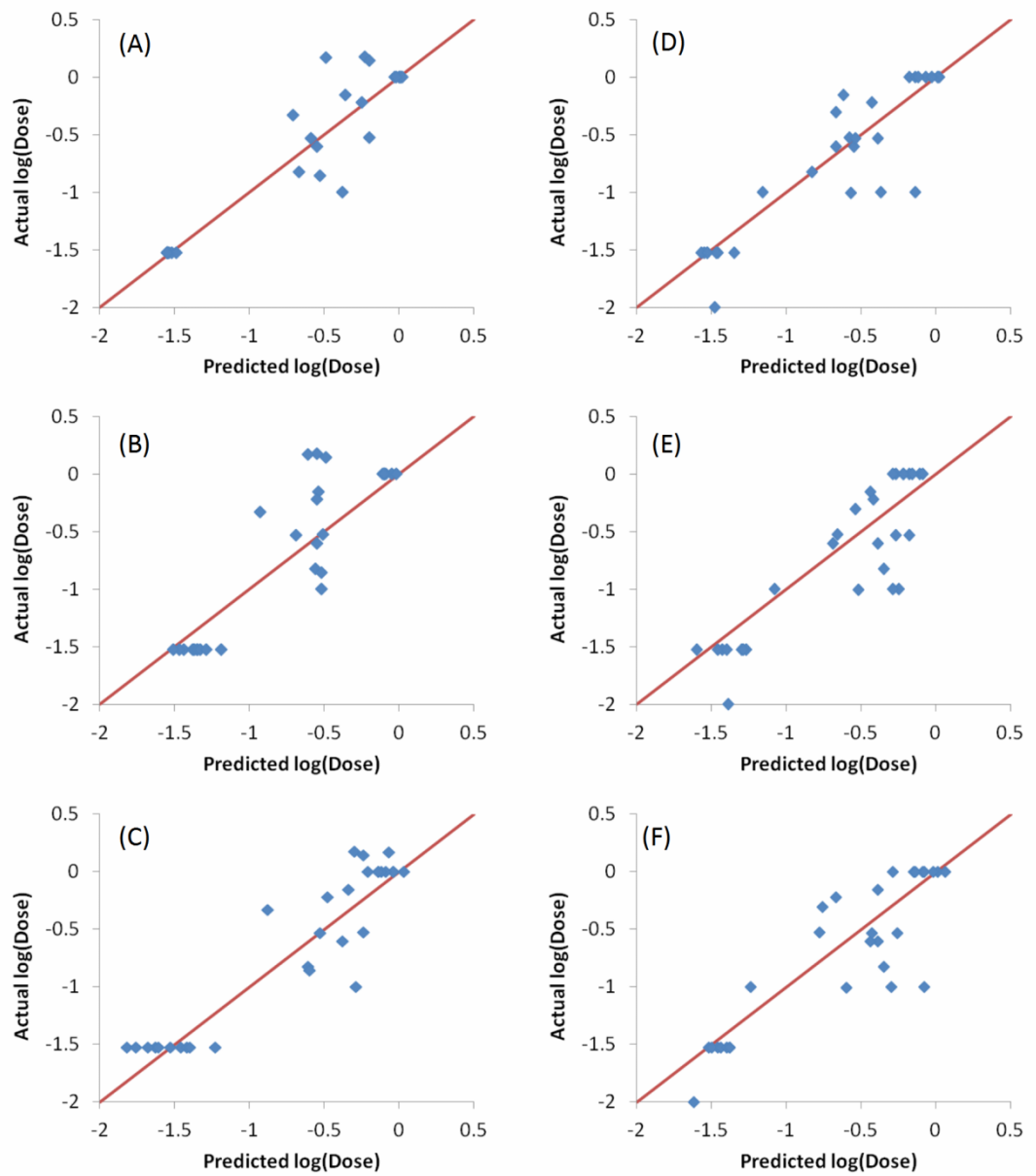
### 3.4. Evaluation of predictive performance by predicting validation set responses

Figure 4 shows that the three machine learning methods resulted in good validation sets predictions, as evident from the predicted points being close to the straight lines (shown in red and representing perfect correlation) in the actual vs predicted plots. It is also clear that the different machine learning models were capable of differentiating between the lipids (entries) with low log(dose), which are the desirable lipids (or formulations), and the lipids/formulations with higher doses.

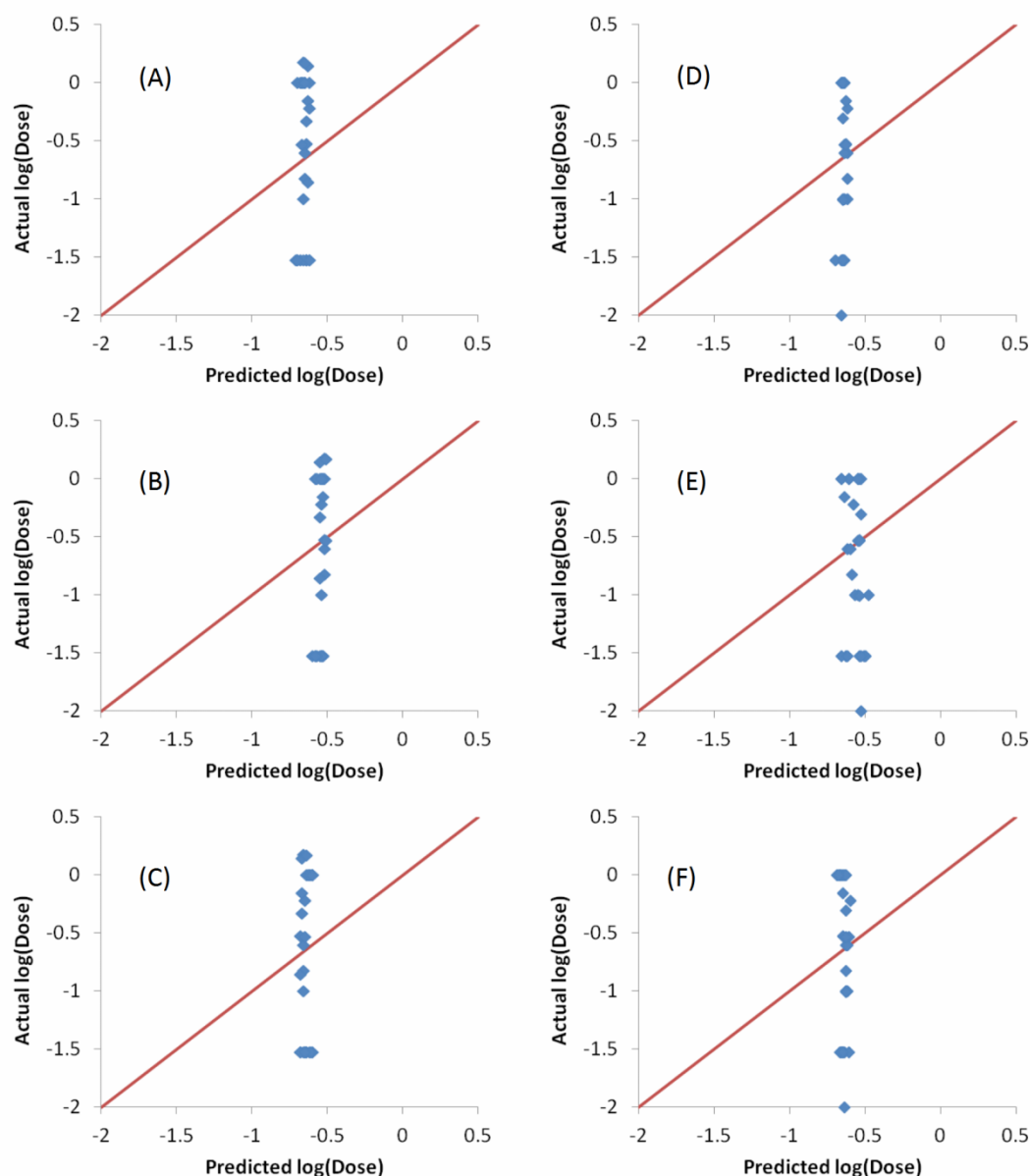
**Table 5.** Evaluation of predictive performance of the different machine learning models.

Set	Machine learning Model	RMSE <sub>val</sub>	R <sup>2</sup> <sub>val</sub>
1	ANN	0.23	0.89
1	SVM	0.32	0.81
1	PLS	0.26	0.86
2	ANN	0.30	0.80
2	SVM	0.36	0.72
2	PLS	0.34	0.75

The curated scaled data set together with an example of the resulting 300 predictor combinations (training set 1) after selection by the evolutionary algorithm and ANN is provided as supplementary materials. An R script for calculating the median predictions of validation set 1 and the associated R<sup>2</sup><sub>val</sub> and RMSE<sub>val</sub> using the data set and the descriptor combinations is also provided as supplementary material.



**Figure 4.** Actual vs predicted log(dose) plots. A-C: Validation set 1, A: ANN, B: SVM and C: PLS. D-F: Validation set 2, D: ANN, E: SVM and F: PLS.



**Figure 5.** Actual vs. predicted responses of validation sets after Y-randomization of training sets responses. A – C: validation set one. A: ANN, B: SVM and C: PLS. D – F: validation set two. D: ANN, E: SVM and F: PLS.

### 3.5. Y-randomization of training set responses

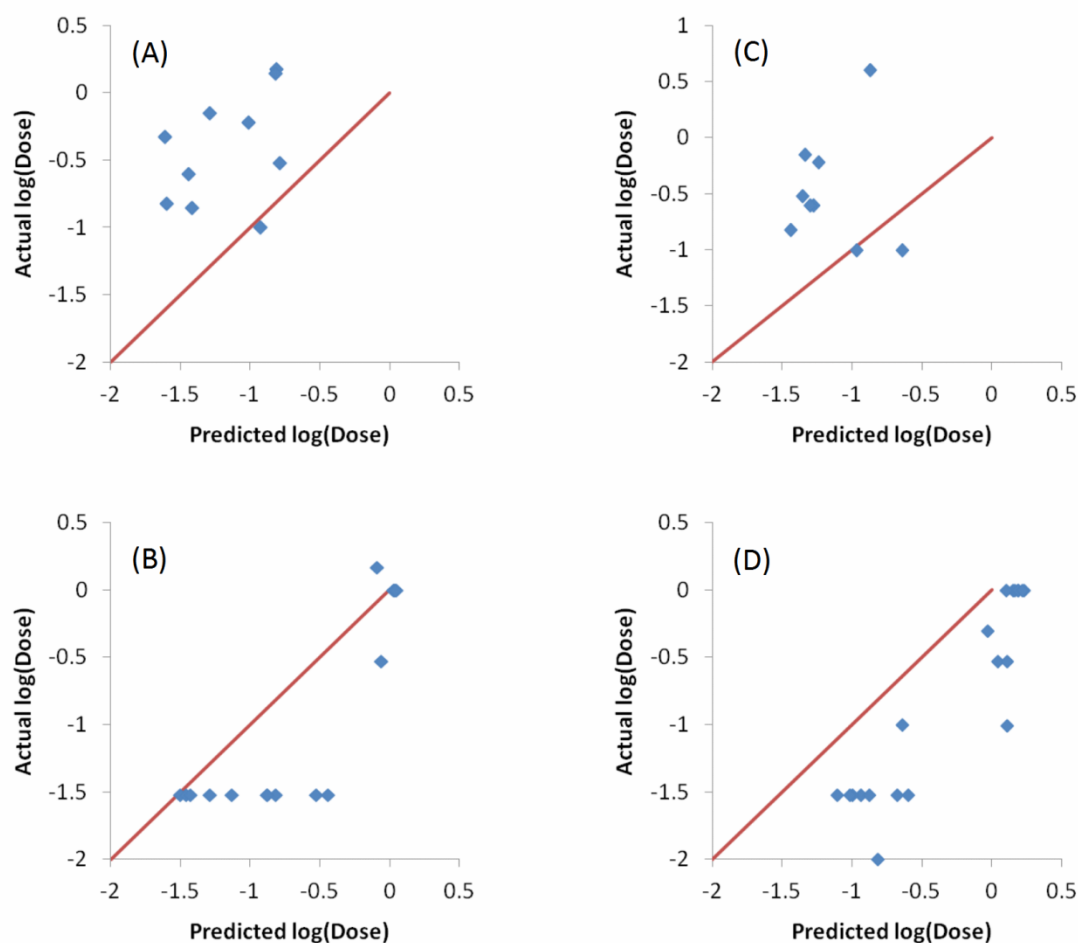
Y-randomization involves randomizing the responses column and then training the predictive models using one of the machine learning methods, with the input descriptors and the responses being mismatched due to the randomization of the responses [46]. Y-randomization was carried-out using the final combinations selected by the evolutionary algorithm as inputs. The resulting predictions together with the actual responses are shown in Figure 5. It can be seen that there is no correlation between the predicted and actual responses for both validation sets and for all of the machine learning methods used. The  $R^2_{val}$  values ranged from 0.014 to 0.116, with  $RMSE_{val}$  values between 0.66 and 0.68. This lack of correlation proves that the results obtained without randomization of the responses (Figure 4 and Table 5) were not due to random chance.

### 3.6. Effect of setting the formulation descriptor PEG mol % to either the maximum or the minimum value

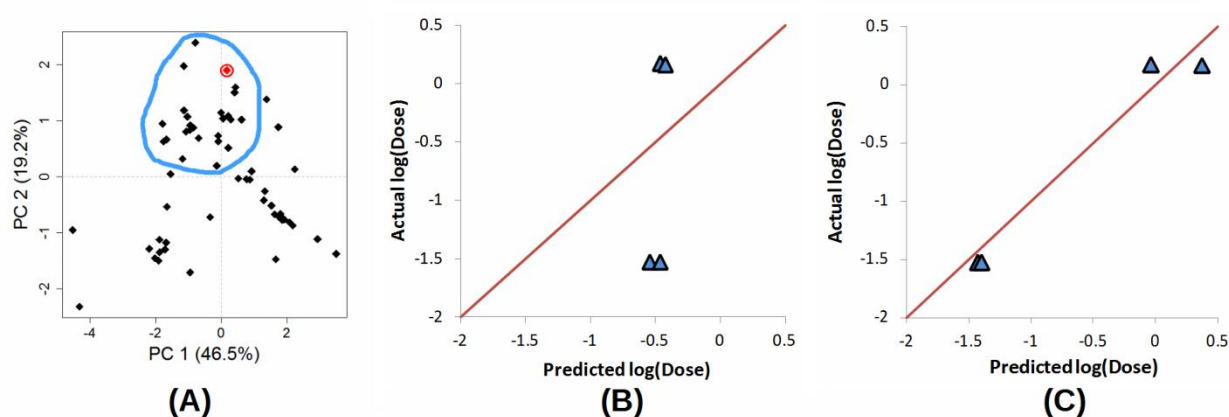
To examine if the predictive models capture the changes in the formulation descriptor; the PEG mol %, the values of this descriptor were set to either its maximum value or rather its minimum counterpart. It is well known that when using siRNA lipoplexes, there is a certain PEG mol % that results in the maximum *in vivo* efficacy in addition to stabilization of the nanoparticles [36,47,48]. The general trend is that increasing the PEG mol % more than a specific mole percent results in decreasing the *in vivo* efficacy. It is generally found that PEG mol % that is equal to 10 decreases efficacy, while values around 1.5% results in good *in vivo* efficacy [1,36]. Hypothetically, it is assumed that if the PEG mol % descriptor values were set to the maximum (equivalent to 10%), the *in vivo* efficacy should decrease, i.e., the log(dose) should increase. On the other hand, if the PEG mol % values are set to the minimum (equivalent to 1.5%), then the *in vivo* efficacy should generally improve for the validation sets lipids that have PEG mol % higher than 1.5%.

It can be seen in Figure 6A and 6C that setting the PEG mol % to the minimum values resulted in a decrease in log(dose) as expected, as evident by the shift of the predictions towards the left hand side. Similarly, setting PEG mol % to the maximum value resulted in shifting of the predicted log(dose) towards higher values as it would be expected (Figure 6B and 6D). These results prove that the predictive models were capable of capturing the significance of the formulation descriptor in a correct manner. ANN was the method used to train the models because it resulted in the best predictions as shown in Figure 4 and Table 5. Similar results were obtained with SVM and PLS regression (data not shown).





**Figure 6.** Actual vs. predicted responses of validation sets after setting the values of the PEG mol % descriptor to either the minimum value (A and C) or the maximum value (B and D). The modeling was carried out by ANN. A and B: validation set one. C and D: validation set 2. The validation sets entries with the actual PEG mol% being the maximum value were omitted from A and C, while those with the actual PEG mol % being the minimum were omitted from B and D for visualization clarity.



**Figure 7.** Determination of the applicability domain (AD) of four lipids from validation set one. A: PCA of training set together with one of the validation set lipids (lipid 15) shown in red circle. B: The actual vs predicted plot before determining AD. C: The actual vs predicted plot after determining AD. Predictions in C and B are carried out by ANN. The red line in B and C represents perfect correlation between actual and predicted values.

### 3.7. Refining the predictions by determining the applicability domain (AD)

AD represents a theoretical region in the chemical space of the training set samples. It is expected that predicting the response of unknown samples, e.g., an external validation set, results in more reliable predictions when the unknown samples falls within this region [49,50]. One method to determine this region is by applying PCA on the training and validation data, and constructing the region of applicability accordingly [49]. Figure 7A shows the score plot of one fold of training set one and lipid 15 which belongs to validation set one (shown as a red circle). The descriptors combination used to perform PCA were chosen randomly from one of the final combinations selected by the evolutionary algorithm. The region encircled by the blue line is the AD, and it was determined manually by excluding from the training entries under consideration those which are far from lipid 15 in the space generated by plotting PC 1 and PC 2. The first two components capture 66% of the variance in the data. The training lipids selected within the AD were then used by ANN to predict the response of lipid 15. This procedure was repeated for another three lipids from the same validation set. The four lipids selected were chosen based on them having the highest biases in their predicted values (Table 6). It is clear by comparing the predicted responses in Table 6 before and after carrying out the selection of training lipids lying in the AD that there was a vast improvement in the quality of the predictions as seen from the much lower bias values before and after selection. In addition, the  $R^2$  for the four lipids was 0.47 and 0.96 before and after applying AD lipid selection respectively, showing significant improvement in the prediction accuracy of these lipids. The impact of improvement of predictions can be seen in Figure 7 B and C, where the predictions lies much closer to the red line in Figure 7C compared to 7B. Since this procedure is carried out manually, we suggest that it should be performed as a refining step for the set of lipids that will be chosen for further wet lab experimentations.

**Table 6.** Refinement of predictions by selecting training lipids within AD

Lipid index	Actual response Log(dose)	Predicted response before applying AD	Predicted response after applying AD	Bias before applying AD selection	Bias after applying AD selection
15	-1.52	-0.55	-1.43	0.97	0.09
16	-1.52	-0.47	-1.40	1.05	0.12
70	0.18	-0.47	-0.04	-0.65	-0.22
109	0.17	-0.43	0.37	-0.60	0.20

#### 4. Discussion

This study provides a computational framework to predict *in silico* the *in vivo* performance of the siRNA lipid nanoparticles. The main question answered in this manuscript is how to predict the siRNA dose of siRNA lipid nanoparticles given a set of molecular descriptors, formulation characteristics and a required knockdown percent. From the results presented in this work, it is evident that this objective was successfully achieved. In order to produce high quality predictions, the following aspects were carefully considered; (1) The selection of the optimal descriptor combinations (2) The modeling approach (3) Validation of the machine learning models using external validation sets and (4) Improving the predictive outcome of the final models by selecting the training set lipids according to the applicability domain.

When preparing the data set, 2D descriptors were calculated from the ionizable lipid structures rather than 3D descriptors. The reason for avoiding the use of 3D descriptors is that not all the lipids were defined in terms of their stereochemistry. In addition, the optimized 3D structure of a single molecule present in the solution state might differ from the 3D structure of the same molecule if present in close contact with other molecules as

in the case of nanoparticles. The effect of the source of the 3D structure and its preparation method and energy minimization in relation to the quality of predictions of three classes of molecules (anilines, carboxylic acids and phenols) has been previously shown [51]. There are other potentially important formulation factors that may play a role in the modeling, e.g., particle size and siRNA to lipid ratio, however, they were not included as they were not reported consistently in the selected literature. For example, particle size was reported on occasions as a wide range instead of well defined values. Nanoparticles  $pK_a$  was also not included in the descriptors as it is not initially a controllable variable that could be pre-determined compared to the formulation parameters, the lipid structure (by its design) and the required percent knock-down.

As for the descriptor selection, an evolutionary algorithm was used. The evolutionary algorithm comprised: (a) 'selection' of the descriptor combinations based on an optimization criterion; the RMSE of the test set after splitting the training set into three folds during training, (b) 'crossover' of the selected parent combinations to make new offspring combinations and (c) 'mutations' of certain descriptors in offspring combinations. These processes are main elements in any evolutionary algorithm [52]. Evolutionary algorithms are suitable for solving the problem of finding optimized solutions of combinations from a set of inputs (descriptors in this case) where an exhaustive search that covers all possible combinations is computationally not feasible [53]. Accordingly, evolutionary algorithms and their variants, such as genetic algorithms, were used to refine the structure of Au nanoparticles [54] and to optimize descriptor combinations in counter-propagation artificial neural networks models used to classify drugs as being either hepatotoxic or nonhepatotoxic [55].

The modeling approach in the current work involved three machine learning methods; ANN, SVM and PLS. These methods differ in their inner workings. The ANNs are considered a collection of linear and non-linear functions that are governed by the choice of the ANN architecture and activation functions. The SVM belongs to the class of kernel algorithms while PLS regression depends on the construction of latent components (principle components) that result in the best covariance with the response variable. Thus, the difference in their predictive performance could be expected. In order to improve the predictive outcome of the final models, averaging of the predicted response values was carried out. Averaging of predictions belongs to a set of machine learning methods called ensemble learning, and usually results in better prediction outcome [56].

Machine learning models require reliable validation to be sure about their ability to successfully predict unknown observations responses. For this purpose, many metrics were suggested and used such as  $R^2$ ,  $Q^2$  and external validation set  $R^2$ . Similarly, RMSE of training set predictions, cross-validation RMSE and external validation RMSE are used for the same purpose. In addition, techniques such as Y-randomization are used to exclude the possibility of the model predictions being due to random chance.  $Q^2$ , the cross-validation coefficient of determination, does not necessarily correlate with good predictive performance for external validation sets [57]. Thus, in this work the validation of the final machine learning models was carried out by predicting responses of two external validation sets as well as performing Y-randomization of training set responses, conforming to the best model validation practices [50,58]. The results showed that the obtained models are reliable.

It is suggested that training set composition and/or the relevant properties of the validation set in relation to the training set governs the predictive performance [59,60]. One way to overcome this is to make sure that the validation set observations are within the applicability domain of the training set [50,58]. In the current work, rather than selecting the validation set observations that lie within the training set applicability domain, a reverse approach was followed; a subset of the training set elements were selected to be close in the predictor space to the validation element under investigation, i.e., these selected training set elements were used to construct the applicability domain. PCA of the training set and the validation set lipid was carried out to determine this applicability domain visually (Figure 7A). It is evident from the results presented in Figure 7B

and 7C and Table 6 that this protocol resulted in significant improvement in performance.

Recently, *in vitro* cellular uptake of siRNA nanoparticles formulated with hydrophobic derivatives of polyethyleneimine (PEI) was predicted by QSAR modeling using either linear regression, random forests or multilayer perceptron, with the nonlinear methods proving to be more efficient than linear regression [61]. The  $R^2$  of the external test set ranged between 0.34 to 0.50 depending on the machine learning method used and on the number of input descriptors, with the initial number of 26 descriptors being reduced either by binary encoding or by backward elimination.

Overall, in the current work for the first time, *in vivo* performance of siRNA nanoparticles could be predicted accurately by combining machine learning techniques with cheminformatics. This framework will greatly enhance the development of siRNA nanomedicines.

## 5. Conclusions

The *in vivo* efficacy of siRNA ionizable lipid nanoparticles could be predicted with excellent accuracy provided careful modeling choices. Calculating molecular descriptors of a series of ionizable lipids followed by selecting best descriptor combinations using an evolutionary algorithm in combination with machine learning modeling by ANN, SVM and PLS and then finally making an ensemble of the predictions by calculating the median of validation set predictions resulted in successful predictions of *in vivo* activity of siRNA ionizable lipids nanoparticles. Depending on the machine learning method and the validation set,  $R^2_{val}$  of up to 0.89 could be achieved. Further improvement of validation set entries with high bias was achievable by selecting training lipids within the applicability domain, with  $R^2_{val}$  improvement from 0.47 to 0.96.

This is the first study to predict *in vivo* performance of siRNA lipoplexes formulated with ionizable lipids, based on the lipids structure and certain nanoparticle characteristics. This *in silico* approach allows the evaluation of virtually an endless number of ionizable lipids prior to their actual synthesis and wet lab evaluation and hence saving valuable resources and time while exploring the vast chemical space of these lipids and their formulations.

**Supplementary Materials:** The following are available online at [www.mdpi.com/xxx/s1](http://www.mdpi.com/xxx/s1), Data set: data\_set.csv, example of descriptor combinations for validation set 1 using ANN: combinations.RData, R script to calculate the median of predictions using the provided descriptor combinations: predict\_dose.R.

**Author Contributions:** Conceptualization, A.A.M.; methodology A.A.M, A.A.N and R.M.H; data curation and software, A.A.M; formal analysis, A.A.M., A.A.N and R.M.H; project administration, A.A.M.; writing—original draft preparation, A.A.M., A.A.N and R.M.H; writing—review and editing, A.A.M., A.A.N and R.M.H. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research received no external funding.

**Data Availability Statement:** Data is contained within the article or supplementary material.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

1. Jayaraman, M.; Ansell, S.M.; Mui, B.L.; Tam, Y.K.; Chen, J.; Du, X.; Butler, D.; Eltepu, L.; Matsuda, S.; Narayanannair, J.K.; et al. Maximizing the Potency of siRNA Lipid Nanoparticles for Hepatic Gene Silencing In Vivo. *Angew. Chem. Int. Ed.* **2012**, *51*, 8529–8533, doi:10.1002/anie.201203263.
2. Sato, Y.; Hashiba, K.; Sasaki, K.; Maeki, M.; Tokeshi, M.; Harashima, H. Understanding structure-activity relationships of pH-sensitive cationic lipids facilitates the rational identification of promising lipid nanoparticles for delivering siRNAs in vivo. *J. Control. Release* **2019**, *295*, 140–152, doi:https://doi.org/10.1016/j.jconrel.2019.01.001.
3. Molla, M.R.; Chakraborty, S.; Munoz-Sagredo, L.; Drechsler, M.; Orian-Rousseau, V.; Levkin, P.A. Combinatorial Synthesis of a Lipidoid Library by Thiolactone Chemistry: In Vitro Screening and In Vivo Validation for siRNA Delivery. *Bioconj. Chem.* **2020**, *31*, 852–860, doi:10.1021/acs.bioconjchem.0c00013.

4. Zhang, M.M.; Bahal, R.; Rasmussen, T.P.; Manautou, J.E.; Zhong, X.-b. The growth of siRNA-based therapeutics: Updated clinical studies. *Biochem. Pharmacol.* **2021**, *189*, 114432, doi:https://doi.org/10.1016/j.bcp.2021.114432.
5. Fire, A.; Xu, S.Q.; Montgomery, M.K.; Kostas, S.A.; Driver, S.E.; Mello, C.C. Potent and specific genetic interference by double-stranded RNA in *Caenorhabditis elegans*. *Nature* **1998**, *391*, 806-811.
6. Elbashir, S.M.; Harborth, J.; Lendeckel, W.; Yalcin, A.; Weber, K.; Tuschl, T. Duplexes of 21-nucleotide RNAs mediate RNA interference in cultured mammalian cells. *Nature* **2001**, *411*, 494-498.
7. Titze-de-Almeida, R.; David, C.; Titze-de-Almeida, S.S. The Race of 10 Synthetic RNAi-Based Drugs to the Pharmaceutical Market. *Pharm. Res.* **2017**, *34*, 1339-1363, doi:10.1007/s11095-017-2134-2.
8. Dong, Y.; Siegwart, D.J.; Anderson, D.G. Strategies, design, and chemistry in siRNA delivery systems. *Adv. Drug Del. Rev.* **2019**, *144*, 133-147, doi:https://doi.org/10.1016/j.addr.2019.05.004.
9. ClinicalTrials.gov. Available online: [https://clinicaltrials.gov/ct2/results?term=siRNA&age\\_v=&gndr=&type=&rslt=&phase=1&phase=2&phase=3&Search=Apply](https://clinicaltrials.gov/ct2/results?term=siRNA&age_v=&gndr=&type=&rslt=&phase=1&phase=2&phase=3&Search=Apply) (accessed on 17/03/2020).
10. Dowdy, S.F. Overcoming cellular barriers for RNA therapeutics. *Nat. Biotechnol.* **2017**, *35*, 222-229, doi:10.1038/nbt.3802.
11. Whitehead, K.A.; Langer, R.; Anderson, D.G. Knocking down barriers: Advances in siRNA delivery. *Nat. Rev. Drug Discov.* **2009**, *8*, 129-138, doi:10.1038/nrd2742.
12. Paunovska, K.; Gil, C.J.; Lokugamage, M.P.; Sago, C.D.; Sato, M.; Lando, G.N.; Gamboa Castro, M.; Bryksin, A.V.; Dahlman, J.E. Analyzing 2000 in Vivo Drug Delivery Data Points Reveals Cholesterol Structure Impacts Nanoparticle Delivery. *ACS Nano* **2018**, *12*, 8341-8349, doi:10.1021/acsnano.8b03640.
13. Cullis, P.R.; Hope, M.J. Lipid Nanoparticle Systems for Enabling Gene Therapies. *Mol. Ther.* **2017**, *25*, 1467-1475, doi:10.1016/j.yymthe.2017.03.013.
14. Metwally, A.A.; Blagbrough, I.S.; Mantell, J.M. Quantitative silencing of EGFP reporter gene by self-assembled siRNA lipoplexes of LinOS and cholesterol. *Mol. Pharmaceutics* **2012**, *9*, 3384-3395, doi:10.1021/mp300435x.
15. Metwally, A.A.; Reelfs, O.; Pourzand, C.; Blagbrough, I.S. Efficient silencing of EGFP reporter gene with siRNA delivered by asymmetrical N<sup>4</sup>,N<sup>9</sup>-diacyl spermines. *Mol. Pharmaceutics* **2012**, *9*, 1862-1876, doi:10.1021/mp200429n.
16. Wesolowski, M.; Suchacz, B. Artificial Neural Networks: Theoretical Background and Pharmaceutical Applications: A Review. *J. AOAC Int.* **2019**, *95*, 652-668, doi:10.5740/jaoacint.SGE\_Wesolowski\_ANN.
17. Maltarollo, V.G.; Kronenberger, T.; Espinoza, G.Z.; Oliveira, P.R.; Honório, K.M. Advances with support vector machines for novel drug discovery. *Expert. Opin. Drug Discov.* **2019**, *14*, 23-33, doi:10.1080/17460441.2019.1549033.
18. Heikamp, K.; Bajorath, J. Support vector machines for drug discovery. *Expert. Opin. Drug Discov.* **2014**, *9*, 93-104, doi:10.1517/17460441.2014.866943.
19. Hathout, R.M.; Metwally, A.A.; Woodman, T.J.; Hardy, J.G. Prediction of Drug Loading in the Gelatin Matrix Using Computational Methods. *ACS Omega* **2020**, *5*, 1549-1556, doi:10.1021/acsomega.9b03487.
20. Boulesteix, A.L.; Strimmer, K. Partial least squares: a versatile tool for the analysis of high-dimensional genomic data. *Brief Bioinform* **2007**, *8*, 32-44, doi:10.1093/bib/bbl016 [pii]
21. Hathout, R.M.; El-Ahmady, S.H.; Metwally, A.A. Curcumin or bisdemethoxycurcumin for nose-to-brain treatment of Alzheimer disease? A bio/chemo-informatics case study. *Nat. Prod. Res.* **2017**, *Accepted manuscript*.
22. Hathout, R.M.; Abdelhamid, S.G.; El-Housseiny, G.S.; Metwally, A.A. Comparing cefotaxime and ceftriaxone in combating meningitis through nose-to-brain delivery using bio/chemoinformatics tools. *Sci. Rep.* **2020**, *10*, 21250, doi:10.1038/s41598-020-78327-w
23. Hathout, R.M.; Abdelhamid, S.G.; El-Housseiny, G.S.; Metwally, A.A. Comparing cefotaxime and ceftriaxone in combating meningitis through nose-to-brain delivery using bio/chemoinformatics tools. *Sci. Rep.* **2020**, *10*, 21250, doi:10.1038/s41598-020-78327-w [pii].
24. Todeschini, R.; Consonni, V. *Handbook of Molecular Descriptors*; WILEY-VCH Verlag GmbH: 2008.
25. Todeschini, R.; Consonni, V.; Todeschini, R. Molecular Descriptors. In *Recent Advances in QSAR Studies: Methods and Applications*, Leszczynski, J., Puzyn, T., Cronin, M.T.D., Eds.; Challenges and Advances in Computational Chemistry and Physics; Springer: 2010; Volume 8, pp. 29-102.
26. Valdés-Martín, J.R.; Marrero-Ponce, Y.; García-Jacas, C.R.; Martínez-Mayorga, K.; Barigye, S.J.; Vaz d'Almeida, Y.S.; Pham-The, H.; Pérez-Giménez, F.; Morell, C.A. QuBiLS-MAS, open source multi-platform software for atom- and bond-based topological (2D) and chiral (2.5D) algebraic molecular descriptors computations. *J. Cheminformatics* **2017**, *9*, 35, doi:10.1186/s13321-017-0211-5.
27. Shamay, Y.; Shah, J.; Işık, M.; Mizrachi, A.; Leibold, J.; Tschaharganeh, D.F.; Roxbury, D.; Budhathoki-Uprety, J.; Nawaly, K.; Sugarman, J.L.; et al. Quantitative self-assembly prediction yields targeted nanomedicines. *Nat. Mater.* **2018**, *17*, 361-368, doi:10.1038/s41563-017-0007-z.
28. De Benedetti, P.G.; Fanelli, F. Computational modeling approaches to quantitative structure-binding kinetics relationships in drug discovery. *Drug Discov. Today* **2018**, doi:https://doi.org/10.1016/j.drudis.2018.03.010.
29. Kausar, S.; Falcao, A.O. An automated framework for QSAR model building. *J. Cheminformatics* **2018**, *10*, 1, doi:10.1186/s13321-017-0256-5
30. Kausar, S.; Falcao, A.O. An automated framework for QSAR model building. *J. Cheminformatics* **2018**, *10*, 1, doi:10.1186/s13321-017-0256-5
31. Kausar, S.; Falcao, A.O. An automated framework for QSAR model building. *J. Cheminformatics* **2018**, *10*, 1, doi:10.1186/s13321-017-0256-5 [pii].



32. 29. Reker, D.; Perna, A.M.; Rodrigues, T.; Schneider, P.; Reutlinger, M.; Mönch, B.; Koeberle, A.; Lamers, C.; Gabler, M.; Steinmetz, H.; et al. Revealing the macromolecular targets of complex natural products. *Nat. Chem.* **2014**, *6*, 1072–1078, doi:10.1038/nchem.2095
33. <https://www.nature.com/articles/nchem.2095#supplementary-information>.
34. 30. Metwally, A.A.; Hathout, R.M. Computer-Assisted Drug Formulation Design: Novel Approach in Drug Delivery. *Mol. Pharmaceutics* **2015**, *12*, 2800–2810, doi:10.1021/mp500740d.
35. 31. Wang, W.; Sedykh, A.; Sun, H.; Zhao, L.; Russo, D.P.; Zhou, H.; Yan, B.; Zhu, H. Predicting Nano–Bio Interactions by Integrating Nanoparticle Libraries and Quantitative Nanostructure Activity Relationship Modeling. *ACS Nano* **2017**, *11*, 12641–12649, doi:10.1021/acsnano.7b07093.
36. 32. Basant, N.; Gupta, S. Modeling uptake of nanoparticles in multiple human cells using structure–activity relationships and intercellular uptake correlations. *Nanotoxicology* **2017**, *11*, 20–30, doi:10.1080/17435390.2016.1257075.
37. 33. Liu, R.; Rallo, R.; Bilal, M.; Cohen, Y. Quantitative structure-activity relationships for cellular uptake of surface-modified nanoparticles. *Combinatorial Chem. High Throughput Screening* **2015**, *18*, 365–375, doi:CCHTS-EPUB-65712 [pii]
38. 10.2174/1386207318666150306105525.
39. 34. Rajappan, K.; Tanis, S.P.; Mukthavaram, R.; Roberts, S.; Nguyen, M.; Tachikawa, K.; Sagi, A.; Sablad, M.; Limphong, P.; Leu, A.; et al. Property-Driven Design and Development of Lipids for Efficient Delivery of siRNA. *J. Med. Chem.* **2020**, *63*, 12992–13012, doi:10.1021/acs.jmedchem.0c01407.
40. 35. Whitehead, K.A.; Dorkin, J.R.; Vegas, A.J.; Chang, P.H.; Veiseh, O.; Matthews, J.; Fenton, O.S.; Zhang, Y.; Olejnik, K.T.; Yesilyurt, V.; et al. Degradable lipid nanoparticles with predictable in vivo siRNA delivery activity. *Nat. Commun.* **2014**, *5*, 4277, doi:ncomms5277 [pii]
41. 10.1038/ncomms5277.
42. 36. Kumar, V.; Qin, J.; Jiang, Y.; Duncan, R.G.; Brigham, B.; Fishman, S.; Nair, J.K.; Akinc, A.; Barros, S.A.; Kasperkovitz, P.V. Shielding of Lipid Nanoparticles for siRNA Delivery: Impact on Physicochemical Properties, Cytokine Induction, and Efficacy. *Mol. Ther. Nucleic Acids* **2014**, *3*, e210, doi:mtna201461 [pii]
43. 10.1038/mtna.2014.61.
44. 37. Alabi, C.A.; Love, K.T.; Sahay, G.; Yin, H.; Luly, K.M.; Langer, R.; Anderson, D.G. Multiparametric approach for the evaluation of lipid nanoparticles for siRNA delivery. *Proc. Natl. Acad. Sci. U.S.A.* **2013**, *110*, 12881–12886, doi:10.1073/pnas.1306529110.
45. 38. O'Boyle, N.M.; Banck, M.; James, C.A.; Morley, C.; Vandermeersch, T.; Hutchison, G.R. Open Babel: An open chemical toolbox. *J. Cheminformatics* **2011**, *3*, 33, doi:10.1186/1758-2946-3-33.
46. 39. Yap, C.W. PaDEL-descriptor: An open source software to calculate molecular descriptors and fingerprints. *J. Comput. Chem.* **2011**, *32*, 1466–1474, doi:https://doi.org/10.1002/jcc.21707.
47. 40. Racz, A.; Bajusz, D.; Heberger, K. Interrelation Limits in Molecular Descriptor Preselection for QSAR/QSPR. *Mol. Inform.* **2019**, *38*, e1800154, doi:10.1002/minf.201800154.
48. 41. Deist, T.M.; Dankers, F.J.W.M.; Valdes, G.; Wijsman, R.; Hsu, I.C.; Oberije, C.; Lustberg, T.; van Soest, J.; Hoebers, F.; Jochems, A.; et al. Machine learning algorithms for outcome prediction in (chemo)radiotherapy: An empirical comparison of classifiers. *Med. Phys.* **2018**, *45*, 3449–3459, doi:10.1002/mp.12967.
49. 42. Hamner, B.; Frasco, M.; LeDell, E. *Metrics: evaluation metrics for machine learning*. , R Package Version 0.1.4; 2018.
50. 43. Karatzoglou, A.; Smola, A.; Hornik, K.; Zeileis, A. kernlab - An S4 Package for Kernel Methods in R. *J. Stat. Softw.* **2004**, *11*, 1–20, doi:10.18637/jss.v011.i09.
51. 44. Mevik, B.-H.; Wehrens, R. The pls Package: Principal Component and Partial Least Squares Regression in R. **2007**, *18*, 23, doi:10.18637/jss.v018.i02.
52. 45. Žuvela, P.; Liu, J.J.; Macur, K.; Bączek, T. Molecular Descriptor Subset Selection in Theoretical Peptide Quantitative Structure–Retention Relationship Model Development Using Nature-Inspired Optimization Algorithms. *Anal. Chem.* **2015**, *87*, 9876–9883, doi:10.1021/acs.analchem.5b02349.
53. 46. Rucker, C.; Rucker, G.; Meringer, M. y-Randomization and its variants in QSPR/QSAR. *J. Chem. Inf. Model.* **2007**, *47*, 2345–2357, doi:10.1021/ci700157b.
54. 47. Mui, B.L.; Tam, Y.K.; Jayaraman, M.; Ansell, S.M.; Du, X.; Tam, Y.Y.C.; Lin, P.J.C.; Chen, S.; Narayanannair, J.K.; Rajeev, K.G.; et al. Influence of Polyethylene Glycol Lipid Desorption Rates on Pharmacokinetics and Pharmacodynamics of siRNA Lipid Nanoparticles. *Mol. Ther. Nucleic Acids* **2013**, *2*, e139, doi:10.1038/mtna.2013.66.
55. 48. Sakurai, Y.; Mizumura, W.; Ito, K.; Iwasaki, K.; Katoh, T.; Goto, Y.; Suga, H.; Harashima, H. Improved Stability of siRNA-Loaded Lipid Nanoparticles Prepared with a PEG-Monoacyl Fatty Acid Facilitates Ligand-Mediated siRNA Delivery. *Mol. Pharmaceutics* **2020**, *17*, 1397–1404, doi:10.1021/acs.molpharmaceut.0c00087.
56. 49. Weaver, S.; Gleeson, M.P. The importance of the domain of applicability in QSAR modeling. *J. Mol. Graphics Model.* **2008**, *26*, 1315–1326, doi:S1093-3263(08)00003-X [pii]
57. 10.1016/j.jmgm.2008.01.002.
58. 50. Tropsha, A. Best Practices for QSAR Model Development, Validation, and Exploitation. *Mol. Inform.* **2010**, *29*, 476–488, doi:10.1002/minf.201000061.



59. 51. Geidl, S.; Svobodová Vařeková, R.; Bendová, V.; Petrusek, L.; Ionescu, C.-M.; Jurka, Z.; Abagyan, R.; Koča, J. How Does the Methodology of 3D Structure Preparation Influence the Quality of pKa Prediction? *J. Chem. Inf. Model.* **2015**, *55*, 1088-1097, doi:10.1021/ci500758w.
60. 52. Sipper, M.; Fu, W.; Ahuja, K.; Moore, J.H. Investigating the parameter space of evolutionary algorithms. *BioData min.* **2018**, *11*, 2, doi:10.1186/s13040-018-0164-x.
61. 53. Douguet, D.; Thoreau, E.; Grassy, G. A genetic algorithm for the automated generation of small organic molecules: Drug design using an evolutionary algorithm. *J. Comput. Aided Mol. Des.* **2000**, *14*, 449-466, doi:10.1023/a:1008108423895.
62. 54. Yu, M.; Yankovich, A.B.; Kaczmarowski, A.; Morgan, D.; Voyles, P.M. Integrated Computational and Experimental Structure Refinement for Nanoparticles. *ACS Nano* **2016**, *10*, 4031-4038, doi:10.1021/acsnano.5b05722.
63. 55. Bajželj, B.; Drgan, V. Hepatotoxicity Modeling Using Counter-Propagation Artificial Neural Networks: Handling an Imbalanced Classification Problem. *Molecules* **2020**, *25*, 481.
64. 56. Oprisiu, I.; Varlamova, E.; Muratov, E.; Artemenko, A.; Marcou, G.; Polishchuk, P.; Kuz'min, V.; Varnek, A. QSPR Approach to Predict Nonadditive Properties of Mixtures. Application to Bubble Point Temperatures of Binary Mixtures of Liquids. *Mol. Inform.* **2012**, *31*, 491-502, doi:10.1002/minf.201200006.
65. 57. Golbraikh, A.; Tropsha, A. Beware of q<sup>2</sup>! *J. Mol. Graphics Model.* **2002**, *20*, 269-276, doi:S1093-3263(01)00123-1 [pii]
66. 10.1016/s1093-3263(01)00123-1.
67. 58. Maleki, F.; Muthukrishnan, N.; Ovens, K.; Reinhold, C.; Forghani, R. Machine Learning Algorithm Validation: From Essentials to Advanced Applications and Implications for Regulatory Certification and Deployment. *Neuroimaging Clin. N. Am.* **2020**, *30*, 433-445, doi:S1052-5149(20)30059-9 [pii]
68. 10.1016/j.nic.2020.08.004.
69. 59. Nalepa, J.; Kawulok, M. Selecting training sets for support vector machines: a review. *Artif. Intell. Rev.* **2019**, *52*, 857-900, doi:10.1007/s10462-017-9611-1.
70. 60. Martin, T.M.; Harten, P.; Young, D.M.; Muratov, E.N.; Golbraikh, A.; Zhu, H.; Tropsha, A. Does Rational Selection of Training and Test Sets Improve the Outcome of QSAR Modeling? *J. Chem. Inf. Model.* **2012**, *52*, 2570-2578, doi:10.1021/ci300338w.
71. 61. Nademi, Y.; Tang, T.; Uludağ, H. Modeling Uptake of Polyethylenimine/Short Interfering RNA Nanoparticles in Breast Cancer Cells Using Machine Learning. *Adv. NanoBiomed Res.* **2021**, 2000106, doi:https://doi.org/10.1002/anbr.202000106.
- 72.